

# METODI ITERATIVI / STAZIONARI

4/12/2012  
(16)

$$Ax = b \quad A \in \mathbb{C}^{n \times n}, x, b \in \mathbb{R}^n \quad x^* = A^{-1}b$$

Si scrive  $A = M - N$  con  $M$  non singolare (nella pratica "facile da invertire")

$$Mx = Nx - b \quad \leadsto \quad x = M^{-1}Nx + M^{-1}b = Bx + c$$

Si definisce  $x^{k+1} = Bx^k + c \quad k \geq 0$  (partendo da un certo  $x_0$  tipicamente 0)

Si dicono metodi STAZIONARI perché la matrice di iterazione  $B$  è sempre la stessa.

Osserviamo che chiaramente vale  $x^* = Bx^* + c$ , inoltre

$$e^k = x^k - x^* \quad \leadsto \quad e^{k+1} = Be^k = B^2e^{k-1} = B^ke^0$$

e dunque il metodo converge (indipendentemente da  $e_0$ ) sse

$$\lim_{k \rightarrow \infty} B^k = 0 \quad \text{sse} \quad \rho(B) < 1$$

OSS: Se posso scomporre  $B = S \Lambda S^{-1}$  (autoval.) la condizione è evidente dato che

$$B^k = S \Lambda^k S^{-1}$$

con  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$  e quindi  $B^k \rightarrow 0$  sse  $(\max_i |\lambda_i|)^k \rightarrow 0$  sse  $\lambda_{\max} < 1$ .

(Nel caso generale non lo dimostra ma vale lo stesso)

**ESEMPLI:** (di metodi) In generale abbiamo:

$$A = D - L - U \quad \text{con} \quad D = \text{diag}(a_{11}, \dots, a_{nn})$$

$$L = - \begin{bmatrix} 0 & & 0 \\ a_{21} & \ddots & \\ & \ddots & 0 \end{bmatrix}$$

$$U = - \begin{bmatrix} 0 & & a_{1n} \\ & \ddots & \\ 0 & & 0 \end{bmatrix}$$

**Jac.** Jacobi:  $x_i^{(k+1)} = \frac{(b_i - \sum_{j \neq i} a_{ij} x_j^{(k)})}{a_{ii}}$

si usa  $M = D$  e  $N = L + U \quad \leadsto \quad B = M^{-1}N = D^{-1}(L + U)$ .

Converge se DDS o DDD + irriducibile.

(Quindi nei nostri casi converge, peccato che sia lentissimo, quindi è meglio!)  
(Pb elitici) CG

**GS.** Gauss-Seidel:  $x_i^{(k+1)} = \frac{(b_i - \sum_{j < i} a_{ij} x_j^{(k+1)} - \sum_{j > i} a_{ij} x_j^{(k)})}{a_{ii}}$

si usa  $M = D - L$  e  $N = U \quad \leadsto \quad B = (D - L)^{-1}U$

Converge per  $A$  sdg ma è un po' lenta!

**SOR.** Si introduce un parametro  $\omega$ :

$$M_\omega = \omega^{-1}D - L$$

$$N_\omega = (1 - \omega)\omega^{-1}D + U$$

dal punto di vista dei conti non cambia, ma così abbiamo da giocare un parametro per aumentare la velocità di convergenza.

Cerchiamo  $\omega^*$  parametro ottimo t.c.

$$\rho(B_{\omega^*}) = \min_{\omega > 0} \rho(B_\omega)$$

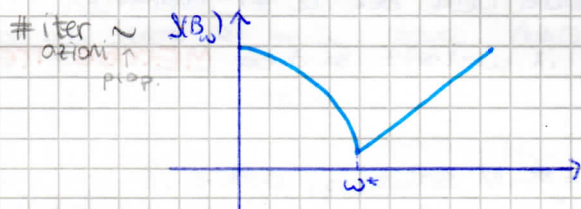


Abbiamo due se  $A$  Sdp tridiagonale o tridiagonale a blocchi allora Jacobi e Gauss-Seidel convergono, GS più velocemente di Jacobi (ma sono anch'essi lenti). Anche SOR converge e si ha

$$w^* = \frac{2}{1 + \sqrt{1 - \rho(B_J)}} \in (1, 2)$$

(che però non è ovvio da calcolare).  $B_J$  matrice di Jacobi

Il grafico di  $\rho(B_w)$  al variare di  $w$  è tipo:



quindi conviene approssimare  $w^*$  per eccesso nel caso non riusciamo a calcolarlo (cost cresce meno velocemente).

In ogni caso anche il SOR perde contro il CG perché la convergenza è simile ma in CG non devo calcolare nessun parametro!

Allora li buttiamo via? No, foranno parte dei metodi MULTIGRID.

### SIMPLE ITERATION

È equivalente al precedente:

$$B = M^{-1}N = I - M^{-1}A$$

$\uparrow$   
 $N = M - A$

Quindi riparto da capo e vedo le cose in un'ottica di preconditionamento.

$Ax = b$  ho scelto un preconditionatore  $M$  t.c.  $M^{-1}A \approx I$  (\*)

$\uparrow$   
in qualche senso che dipende dal metodo

Ci si aspetta che

$$\underbrace{M^{-1}(b - Ax^k)}_{=r^k} \text{ approssimi } e^k = x^* - x^k$$

Infatti:

$$M^{-1}(b - Ax^k) = M^{-1}(Ax^* - Ax^k) = M^{-1}Ae^k \approx e^k \quad (*)$$

$\text{se } (*)$

Viene naturale considerare:

$$x^{k+1} = x^k + M^{-1}(b - Ax^k) \quad (.)$$

infatti se fosse  $M^{-1}A = I$  non in senso approssimato cioè direi meglio se volesse (\*) con l'uguale, avremmo

$$x^{k+1} = x^k + e^k = x^k + x^* - x^k = x^* \quad \text{!}$$

Questo metodo dal punto di vista algebrico è equivalente a quello di prima infatti ho:

$$(*) \rightarrow x^{k+1} = \underbrace{(I - M^{-1}A)}_{=B} x^k + M^{-1}b = M^{-1}N x^k + M^{-1}b \quad (\text{splitting})$$

ALGORITMO (per un qualsiasi metodo di SPLITTING)

Dato  $x^0$ , calcolo  $r^0 = b - Ax^0$  e risolvo  $Mz^0 = r^0$  (almeno non calcolo davvero  $M^{-1}$  !)

per  $k \geq 1$  pongo  $x^k = x^{k-1} + z^{k-1}$

$$r^k = b - Ax^k$$

risolvo  $Mz^k = r^k$   
 $k = k+1$

$\nwarrow$  risolvo esattamente

← da qualche parte c'è un test d'arresto, o subito o prima di risolvere  $Mz^k = r^k$  (che è la cosa pesante)

MA, MI SERVE VERAMENTE? NO  $\rightarrow$  posso usare un metodo iterativo anche qui



Mi rendo conto che posso non calcolare esattamente  $z^{k+1}$ , visto che se sto chiedendo  $x$  con una precisione, non ha senso trovare  $z^k$  con precisione migliore!

In realtà mi interessa essere poco precisa nelle prime iterazioni, ovvero quando sono ancora lontana dalla soluzione "giusta", ma essere più precisa dopo.

Ho un OUTER ITERATION (quella per trovare  $x$ ) con una certa condizione d'arresto (CA) e un INNER ITERATION (quella per trovare  $z^k$ ) con una CA che dipende dalla CA dell'outer iteration.

Se non calcolo esattamente  $z^k$  il metodo si dice **METODO ITERATIVO INESATTO**.

Ora abbiamo il problema, dato  $A=M-N$  splittato, qual'è  $M$  candidato a essere preconditionatore per PCG?

Una condizione che sappiamo deve essere soddisfatta è  $M$  sdp (ovvero se in realtà nella pratica a volte si usa  $M=M'+R$  con  $M'$  sdp e  $R$  "piccola" e se  $R$  non cambia di molto  $M'$  funziona lo stesso).

- Un possibile candidato è quello di "derivare"  $M$  dal metodo di Jacobi:  $M=D=\text{diag}(A)$

Pb ellittico  $\rightarrow A$  sdp  $\Rightarrow a_{ii} > 0$

- SOR simmetrico (SSOR) (GS = SOR per  $w=1$ )

Partiamo da GS, abbiamo:  $M=D-L$ ,  $N=U$

alternativa:  $M=D-U$ ,  $N=L$

(fatto posso partire da  $i=N$  e tornare indietro al posto che da  $i=1$  e andare avanti)

Quindi in generale per il SOR possiamo fare

$$M_w = w^{-1}D - L \quad N_w = (1-w)w^{-1}D + U \quad (\text{SOR})$$

$$\tilde{M}_w = w^{-1}D - U \quad \tilde{N}_w = (1-w)w^{-1}D + L \quad (\text{BACKWARD SOR})$$

L'SSOR consiste nel fare un passo di SOR e uno di backward SOR:

$$\hat{x}^{k+1} = M_w^{-1} (N_w x^k + b)$$

$$x^{k+1} = \tilde{M}_w^{-1} (\tilde{N}_w \hat{x}^{k+1} + b)$$

Nell'implementazione questi due passi vengono mantenuti separati, ma per vedere la matrice di iterazione del metodo possiamo scriverla in un solo passo (facendo dei brutti conti)

$$x^{k+1} = \underbrace{M_w^{-1} N_w M_w^{-1}}_{\tilde{M}^{-1}} N_w x^k + M_w^{-1} (N_w^T + M_w) M_w^{-1} b$$

dove abbiamo considerato il caso  $A=A^T$  e quindi  $M_w^T = \tilde{M}_w$  e  $N_w^T = \tilde{N}_w$ .

Osserviamo che  $\tilde{M}^{-1} = M_w^{-1} N_w^T M_w^{-1}$  è simmetrica (come richiesto).

Si può verificare che:

$$A = \bar{M}_w - \bar{N}_w$$

con  $\bar{M}_w = \frac{w}{2-w} M_w D^{-1} M_w^T$ . (Così vedo che  $\bar{M}_w$  non solo è simmetrica ma è anche sdp, che era l'altro condizione richiesta)

- C'è un ultimo preconditionatore generico (ovvero che non richiede un'analisi degli autovalori del problema e che posso usare per tutti). PRECONDITIONATORE nella FORMA INCOMPLETA di CHOLESKY.



Abbiamo il problema  $Ax=b$ , supponiamo per esempio di essere nel caso delle diff. finite e  $-\Delta u=f$  abbiamo

$$A = \begin{bmatrix} \diagup & & \\ & \diagup & \\ & & \diagup \end{bmatrix} \quad A=LL^T \quad \leadsto \quad \begin{bmatrix} \diagup & & \\ & \diagup & \\ & & \diagup \end{bmatrix} \quad \text{fill in}$$

Si prende  $\tilde{L}$  che è quella "non riempita" ovvero con gli zeri dove  $A$  aveva gli zeri e  $A \neq 0$ . Prendiamo

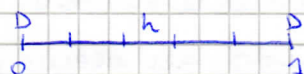
$$M = \tilde{L} \tilde{L}^T \neq A$$

( $\tilde{L}$  si calcola con chol in matlab)  
5/12/2012  
(27)

## MATRICI ASSOCIATE AL PROBLEMA del LAPLACIANO

Consideriamo il problema:

$$\begin{cases} -u'' = f & \text{su } \Omega = (0,1) \\ u(0) = u(1) = 0 \end{cases}$$



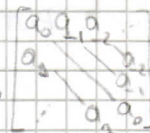
$N = \#$  sottointervalli  
 $N+1$  nodi  
 $h = 1/N$

utilizzando le differenze finite approssimiamo  $-u''(x)$ :

$$-u''(x) \approx \frac{1}{h^2} [-u(x-h) + 2u(x) - u(x+h)]$$

$$\leadsto A_n = \begin{bmatrix} 2 & & & \\ & \ddots & & \\ & & 2 & \\ & & & 2 \end{bmatrix} \quad n = N-1$$

Abbiamo che  $A_n$  è tridiagonale ed è anche una matrice di Toeplitz, ovvero ha elementi costanti lungo le diagonali



Posso interpretare gli  $a_k$  nella matrice di Toeplitz come i coefficienti di Fourier di una funzione  $f$  (**FUNZIONE GENERATRICE di  $A_n$** )

$$a_k = \frac{1}{2\pi} \int_{-\pi}^{\pi} f(t) e^{-itk} dt$$

e quindi ho una corrispondenza tra matrici Toeplitz e funzioni (serie di Fourier) cioè se ho gli  $a_k$  posso trovare  $f$

$$f(t) = \sum_{k \in \mathbb{Z}} a_k e^{itk}$$

Quindi nel nostro caso,  $A_n = \text{tridiag}(-1, 2, -1)$

$$f(t) = 2 - e^{-it} - e^{it} = 2 - 2 \cos t \quad \text{funzione generatrice di } A_n$$

Abbiamo che:

$$\begin{aligned} \lambda_s(A_n) &= f(x_s) \quad \text{con } x_s = \frac{s\pi}{n+1} \quad s=1, \dots, n \\ &= 2 - 2 \cos\left(\frac{s\pi}{n+1}\right) = 4 \sin^2\left(\frac{s\pi}{2(n+1)}\right) \end{aligned}$$

Quindi vediamo che nel nostro caso  $\lambda_1 (= \lambda_{\min})$  va a zero come  $h^2$  (sostituiamo  $n+1 = N = 1/h$ ), mentre  $\lambda_n (= \lambda_{\max})$  rimane limitato ( $\leq 4$ ) e quindi:  $K_2(A_n) \approx h^{-2}$

$A_n$  è anche una  $\tau$  ovvero diagonalizzata da trasformazioni di seni di tipo 1. Quindi

$$A_n = S_n \Lambda S_n$$

con  $S_n$  ortogonale, simmetrica

$$S_n = \frac{\sqrt{2}}{n+1} \left[ \sin\left(\frac{sr\pi}{n+1}\right) \right]_{s,r=1}^n = [S_1 | \dots | S_n] \quad \text{discreta Fourier}$$

ALGEBRA  
 $A = (\lambda_n) A_n = S \Lambda S^T$   
CIRCOLANTI  
 $(\lambda_n) A_n = F_n \Lambda_n F_n^H$   
 $F_n$  trasformata discreta di Fourier  
51



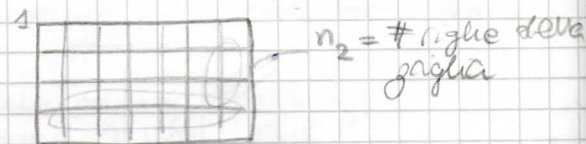
Possiamo ora al caso 2 dimensionale, quindi abbiamo il problema

$$\begin{cases} -\Delta u = f & \text{su } \Omega = (0,1) \times (0,1) \\ u|_{\partial\Omega} = 0 \end{cases}$$

otteniamo la matrice (usando DF)

$$A_n = \begin{bmatrix} 4 & -1 & & \\ -1 & 4 & -1 & \\ & -1 & 4 & -1 \\ & & -1 & 4 \end{bmatrix}$$

in realtà anche usando FEM si ha una triangolazione strutturata uniforme



(ho solo 4 componenti per riga non nulle perché 2 vengano 0 per ortogonalità dei gradienti)

tridiagonale a blocchi, con blocchi tridiagonali. Vediamo qual'è la funzione generatrice, viene:

$$f(u,v) = 4 - 2\cos u - 2\cos v$$

infatti considerando gli elementi di  $A_n$  abbiamo

$$f(u,v) = 4 - e^{i\pi u} e^{i\pi v} - e^{-i\pi u} e^{i\pi v} - e^{i\pi u} e^{-i\pi v} - e^{-i\pi u} e^{-i\pi v}$$

(che è quella sopra).

Posso separare la  $u$  e la  $v$  e scrivere  $f(u,v) = (2 - 2\cos u) + (2 - 2\cos v)$  quindi posso pensare di separare i componenti su  $u$  e  $v$  (per il calcolo degli autovalori).

Autovalori:  $\lambda_{(S_1, S_2)}^{(A_n)} = f\left(\frac{\pi S_1}{n_2+1}, \frac{\pi S_2}{n_2+1}\right)$

$$\begin{aligned} S_1 &= 1, \dots, n_2 \\ S_2 &= 1, \dots, n_2 \end{aligned}$$

(sto immaginando la griglia con notazione doppio indice)

$\lambda_k(A_n)$  con  $k = S_1 + (S_2 - 1)n_2$  se penso all'indicizzazione "torale"

Abbiamo che in questo caso:

$$\lambda_{\max} \leq 8$$

$$\lambda_{\min} \geq h^2$$

$$\rightarrow K_2(A_n) \leq O(h^{-2})$$

$A_n$  è ancora una  $\tau \rightarrow A_n = S_n \Lambda S_n$  con

$$S_n = S_{n_2} \otimes S_{n_2}$$

con  $S_{n_i}$  trasformata discreta di seni di tipo I

def:  $A \in \mathbb{C}^{n \times n}, B \in \mathbb{C}^{m \times m}$

$$A \otimes B = [a_{ij} B]_{m \times m} = \begin{bmatrix} a_{11}B & a_{12}B & a_{13}B \\ a_{21}B & a_{22}B & a_{23}B \end{bmatrix}$$

e quindi:

$$S_{n_2 n_2} = \sqrt{\frac{2}{n_2+1}} \sqrt{\frac{2}{n_2+1}} \left[ \left[ \sin\left(\frac{S_1 \pi t_1}{n_2+1}\right) \right]_{S_1=1, \dots, n_2} \otimes \left[ \sin\left(\frac{S_2 \pi t_2}{n_2+1}\right) \right]_{S_2=1, \dots, n_2} \right]$$

PROP: (PROPRIETÀ PRODOTTO TENSORIALE)  $A, B, C \in \mathbb{C}^{2 \times 2}$  cambio nei vari p.n.

- (i)  $A \otimes (B + C) = A \otimes B + A \otimes C$
- (ii)  $(A + B) \otimes C = A \otimes C + B \otimes C$
- (iii)  $(\alpha A) \otimes B = A \otimes \alpha B = \alpha (A \otimes B) \quad \alpha \in \mathbb{C}$
- (iv)  $A \otimes (B \otimes C) = (A \otimes B) \otimes C$
- (v)  $(A \otimes B)^H = A^H \otimes B^H$

Se  $A \in \mathbb{C}^{n \times n}, B \in \mathbb{C}^{m \times m} \quad (A \otimes B)(C \otimes D) = (AC) \otimes (BD)$

•  $A \otimes B$  non singolare se  $A, B$  non singolari e  $(A \otimes B)^{-1} = A^{-1} \otimes B^{-1}$

•  $\exists \pi \in \mathbb{R}^{nm \times nm}$  di permutazione t.c.  $\pi(A \otimes B)\pi^{-1} = B \otimes A$

• Se  $A \in \mathbb{C}^{n_1 \times n_1}$  e  $B \in \mathbb{C}^{n_2 \times n_2}$  allora:  $\lambda_k(A \otimes B) = \lambda_{i_1}(A) \cdot \lambda_{i_2}(B)$

$$\begin{aligned} k &= 1, \dots, n_1 n_2 \\ i_1 &= 1, \dots, n_1 \\ i_2 &= 1, \dots, n_2 \end{aligned}$$

e gli autovettori sono il prodotto tensoriale di quelli di  $A$  e  $B$ .



•  $\lambda_k(I_{n_2} \otimes A + B \otimes I_{n_1}) = \lambda_{i_1}(A) + \lambda_{i_2}(B)$  con  $A, B$  come prima.

## ESEMPIO: APPLICAZIONE TEOREMA AUTOVALORI di MATRICI HERMITIANE

Vogliamo considerare il problema:

$$\begin{cases} u^{IV} = f & \text{in } (0,1) \\ u^{(r)}(0) = u^{(r)}(1) = 0 & r = 0,1 \end{cases}$$

con le differenze finite ottengo

$$\tilde{A}_n = \begin{bmatrix} 6 & -4 & 1 & 0 \\ -4 & 6 & -4 & 1 \\ 1 & -4 & 6 & -4 \\ 0 & 1 & -4 & 6 \end{bmatrix}$$

che deriva dal fatto che ho:

$$u^{IV}(x) \approx [u(x-2h) - 4u(x-h) + 6u(x) - 4u(x+h) + u(x+2h)] \cdot \frac{1}{h^4}$$

questa formula posso ricavarla in 2 modi, o facendo la serie di Taylor e tutti i conti o considerando:

$$u^{IV} = -(-u'')''$$

sappiamo che  $u''(x) \approx \frac{1}{h^2} [u(x-h) - 2u(x) + u(x+h)]$ , poniamo  $v(x) = -u''(x)$

$$u^{IV} = -v''(x) \approx -\frac{1}{h^2} [v(x-h) - 2v(x) + v(x+h)] =$$

$$= \frac{1}{h^2} [u''(x-h) - 2u''(x) + u''(x+h)]$$

$$= \frac{1}{h^4} [u(x-2h) - 2u(x-h) + u(x) - 2(u(x-h) - 2u(x) + u(x+h)) + \dots + u(x) - 2u(x+h) + u(x+2h)] =$$

$$= \frac{1}{h^4} [u(x-2h) - 4u(x-h) + 6u(x) - 4u(x+h) + u(x+2h)]$$

Mi chiedo, questo trucco funzionerà anche per il calcolo delle funzioni generatrici?

$$-u'' \leadsto f(t) = 2 - 2\cos t$$

$$u^{IV} \leadsto \tilde{f}(t) = 6 - 4e^{it} - 4e^{-it} + e^{i2t} + e^{-2it} = 6 - 8\cos t - 2\cos(2t)$$

Si vede che  $\tilde{f}(t) = f^2(t)$ . Quindi riesco a trovare una relazione tra gli autovalori di  $\tilde{A}_n$  ( $u^{IV}$ ) e di  $A_n$  ( $-u''$ )? (visto che ho calcolato come compimento di  $f$  e  $\tilde{f}$ )

Vediamo la matrice  $\tilde{A}_n$  in questa forma

$$\tilde{A}_n = \begin{bmatrix} 5 & -4 & 1 & & \\ -4 & 6 & -4 & & \\ 1 & -4 & 6 & -4 & \\ & & 1 & -4 & 5 \end{bmatrix} + \begin{bmatrix} 1 & & & & \\ & & & & \\ & & 0 & & \\ & & & 0 & \\ & & & & 1 \end{bmatrix} = A_n^2 + \underline{e_1 e_1^T} + \underline{e_n e_n^T}$$

$\leftarrow$  matrice di correzione, ha rango 2, e sdg  
 $\parallel$  si vede facendo  $A_n^2$  il quadrato di  $A_n$

Vogliamo ricondurre al teorema CF, consideriamo il quoziente di Ray

$$r_{\tilde{A}_n}(x) = \frac{x^T \tilde{A}_n x}{x^T x} = \frac{x^T A_n^2 x}{x^T x} + \frac{x^T e_1 e_1^T x}{x^T x} + \frac{x^T e_n e_n^T x}{x^T x} \geq r_{A_n^2}(x)$$

e quindi per Courant-Fisher  $\lambda_j(\tilde{A}_n) \geq \lambda_j^0(A_n^2) = (2 - 2\cos(\frac{j\pi}{n+1}))^2$



e dunque l'autovalore minimo, che si ottiene per  $J=1$ , è

$$\geq \left( 4 \sec^2 \left( \frac{\pi}{2(n+1)} \right) \right)^2 \approx \frac{\pi^4}{n^4}$$

e dato che il  $\lambda_{\max}$  è limitato si ottiene che

$$K_2(\hat{A}_n) \text{ è al più } O(n^4).$$

Ora ci troviamo nel caso di connessione di  $rk \leq 1$  (Interallacciamento):

$$A \quad B = A + \sigma_1 u_1 u_1^H \quad C = B + \sigma_2 u_2 u_2^H$$

$\nwarrow \hat{A}_n$        $\nwarrow \hat{A}_n$

i teoremi visti ci dicono che

$$\lambda_{J-1}(A) \geq \lambda_J(B) \geq \lambda_J(A) \\ \lambda_{J-1}(B) \geq \lambda_J(C) \geq \lambda_J(B) \quad J=2, \dots, n$$

N.B. in questi teoremi gli autovalori sono indicizzati al contrario

e mettendo insieme le due relazioni ho:

$$\lambda_{J-2}(A) \geq \lambda_{J-1}(B) \geq \lambda_J(C) \geq \lambda_J(B) \geq \lambda_J(A) \quad J=3, \dots, n$$

ovvero ho che controllo tutti gli autovalori di  $C$  ( $\hat{A}_n$ ) con gli autovalori di  $A$  ( $\hat{A}_n^2$ ) tranne  $\lambda_3$  e  $\lambda_2$  (che però in questa notazione sono i più grandi).

In realtà ora mostriamo con un argomento di immersione che possiamo controllare anche  $\lambda_3$  e  $\lambda_2$ .

Considero

$$A_{n+2}^2 = \begin{bmatrix} 5 & -4 & 1 \\ -4 & \hat{A}_n & 1 \\ 1 & -4 & 5 \end{bmatrix}$$

ovvero immerso  $\hat{A}_n$  in  $A_{n+2}^2$ .

sono le dimensioni delle matrici

$$\begin{array}{ccc} A_{n+2} & B_{n+1} & C_n \\ \uparrow & & \uparrow \\ A_{n+2} & & \hat{A}_n \end{array}$$

$$\text{Teo di proiezione} \rightarrow \lambda_J(A) \geq \lambda_J(B) \geq \lambda_{J+1}(A) \quad J=1, \dots, m-1 \quad m=n+2 \\ \lambda_J(B) \geq \lambda_J(C) \geq \lambda_{J+1}(B) \quad J=1, \dots, m-2$$

componendole come prima:

$$\lambda_J(A) \geq \lambda_J(C) \geq \lambda_{J+1}(B) \geq \lambda_{J+2}(A) \quad J=1, \dots, \underline{m-2} = n$$

e quindi controllo tutti gli autovalori di  $C$  ( $= \hat{A}_n$ ).

11/12/2012

## METODI MULTIGRID

Partiamo dal problema unidimensionale

$$\begin{cases} -u'' = f & \text{su } (0,1) \\ u(0) = u(1) = 0 \end{cases}$$

Trovate DF o FE otteniamo la matrice  $A_n = \begin{bmatrix} 2 & -1 & & \\ -1 & 2 & -1 & \\ & -1 & 2 & -1 \\ & & -1 & 2 \end{bmatrix}$

Vediamo cosa succede se applichiamo Jacobi alla risoluzione del nostro sistema lineare.

$$\text{Abbiamo } A = D - L - U \quad M = D \quad N = L + U$$

$$\text{Matrice di iterazione: } B_J = D^{-1}(L+U)$$



Introduciamo un metodo di Jacobi pesato:

$$x^{k+1} = B_{J,w} x^k + w \underbrace{D^{-1}b}_{=c}$$

con  $B_{J,w} = (1-w)I + wB_J$ .

Abbiamo

$$\begin{aligned} B_{J,w} &= (1-w)I + wD^{-1}(L+U) = \\ &= (1-w)I + wD^{-1}(D-A) = \\ &= (1-w)I + wI - wD^{-1}A = \\ &= I - wD^{-1}A \end{aligned}$$

Ricordiamo che stiamo usando  $N+1$  nodi equispaziati,  $h = 1/N$ ,  $n = N-1$ .  
Abbiamo che nel nostro caso:

$$B_{J,w} = I - \frac{w}{2} \begin{bmatrix} 2 & & & \\ & -1 & & \\ & & \ddots & \\ & & & -1 & 2 \end{bmatrix} \quad \text{simmetrica.}$$

$$\lambda_i(B_{J,w}) = 1 - \frac{w}{2} \lambda_i(A) \quad \text{e sappiamo che } \lambda_i(A) = 2 - 2\cos\left(\frac{i\pi}{N}\right) \quad i = 1, \dots, n$$

$$= 4 \sin^2\left(\frac{i\pi}{2N}\right)$$

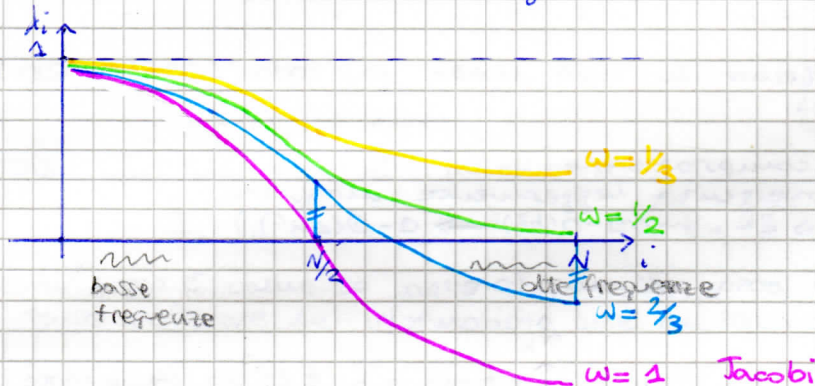
$$w^{(i)} = \left[ \sin\left(\frac{i\pi}{N}\right) \right]_{i=1, \dots, N} \quad \text{autovettori.}$$

Consideriamo l'errore iniziale  $e_0$ , possiamo scriverlo come comb. lin degli autovettori (che sono una base)

$$\Rightarrow e_0 = \sum_{i=1}^{N-1} c_i w^{(i)}$$

Da cui otteniamo:  $e_k = B_{J,w}^k e_0 = \sum_{i=1}^{N-1} c_i B_{J,w}^k w^{(i)} = \sum_{i=1}^{N-1} c_i \lambda_i^k(B_{J,w}) w^{(i)}$

Consideriamo ora l'andamento degli autovalori al variare di  $w$ .



Abbiamo che per  $w \in (0,1)$

$$|\lambda_i(B_{J,w})| < 1 \quad \forall i$$

$$\lambda_1(B_{J,w}) = 1 - 2w \sin^2\left(\frac{\pi}{2N}\right) \underset{N=h}{\approx} 1 - 2w \sin^2\left(\frac{\pi h}{2}\right) = 1 - O(h^2)$$

più grande

Quindi l'autovettore più grande rimane sempre vicino a 1, indipendentemente da  $w$ !  $w$  non ha nessun ruolo da giocare su  $\lambda_1$ .

Inoltre osserviamo che al tendere di  $h$  a zero  $\lambda_1 \rightarrow 1$  e quindi il metodo rallenta.

L'osservazione chiave è che  $w$  può giocare un ruolo in alcune frequenze mentre altre non le tocca.



Per come sono fatti gli autovettori ho che i piccolo corrisponde alle basse frequenze e i grande alle alte frequenze.

Ho che per le basse frequenze  $\omega$  non va mai bene, e non ci posso fare nulla, mentre per le alte frequenze  $\omega$  può giocare un ruolo, perché vediamo che per esempio per  $\omega = 1/2$  gli ultimi autovettori vanno a zero e li possiamo "eliminare".

(Nella pratica si sceglie  $\omega = 2/3$  perché così ho  $|1/\lambda_{N/2}| = |1/\lambda_N|$  e invece di sintonarli a zero so che li limito entrambi e assieme anche quelli compresi).

def: **SMOOTHING FACTOR**:  $s = \max_{\frac{1}{2} \leq i \leq N} |1/\lambda_i(B_{\omega})|$ .

Ho che per  $\omega = 2/3 \rightarrow s = 1/3$   $|1/\lambda_{N/2}| = |1/\lambda_N| = 1/3$  e  $|1/\lambda_i| < 1/3$   $N/2 < i < N$ .

**OSS**: Anche Gauss-Seidel e gli altri metodi stazionari si comportano analogamente, ovvero che non riescono a smorzare le basse frequenze ma riescono a smorzare le alte.

Ora dobbiamo pensare di "spezzare il problema" in più spazi (visto che ho come diviso il mio problema tra basse e alte frequenze e riesco a sistemare le cose da una parte sì e una no)

Al posto di considerare  $Ax=b$  passiamo all'equazione del residuo

$$Ae=r.$$

Osserviamo che i due sistemi sono equivalenti infatti: sia  $x^* = A^{-1}b$

$$e = x^* - \tilde{x} \quad \text{con } \tilde{x} \text{ approssimazione di } x^*$$

$$Ae = A(x^* - \tilde{x}) = b - A\tilde{x} = r$$

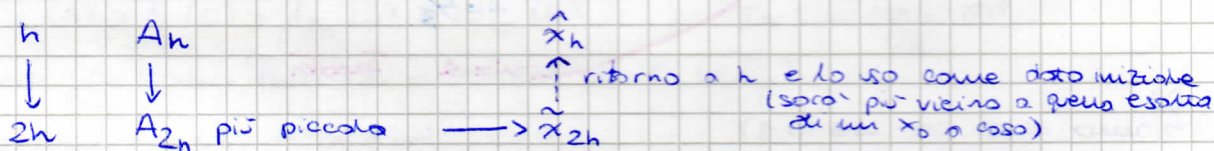
Poi si recupera  $x^*$  facendo  $x^* = \tilde{x} + e$ .

Ora l'idea è di "confondere" il metodo per sistemare anche le basse frequenze.

Proiettiamo il problema su una mesh più rada (COARSE GRID); così abbiamo anche:

- minor costo computazionale
- velocità di convergenza leggermente maggiore (tipo se  $h \rightarrow 2h$  ma  $1-O(h^2) \rightarrow 1-O(4h^2)$ )

Vediamo un'idea di cosa possiamo fare:



Punto chiave:  $Ae=r$ .

Dopo smoother (J, Jw, GS) l'errore è smooth (bassa frequenza)



$x$  = griglia più fissa  
 $\nearrow$  interpolazione lineare

UNA FUNZIONE SMOOTH VIENE APPROSSIMATA BENE CON L'INTERPOLAZIONE LINEARE!  
 (Se non è smooth no!)

Quindi se proietto il problema su una griglia più grande non perdo molto.



Un'altra cosa bella è che il fatto di essere alto o bassa frequenza dipende dallo mesh! Ovvero una cosa che magari è in medio frequenza su uno mesh più fitto può diventare in alto frequenza su uno mesh più rado! (perché le frequenze sono indipendenti dallo mesh, è la loro "classificazione" che cambia.)

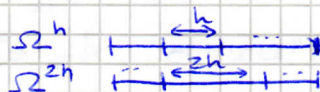
IDEA QUALITATIVA:

**SMOOTHER** • Rilasso  $Ax=b$  su  $\Omega^h$

(applico lo smoother  $B_{s,w}, B_{gs}, \dots$ )

ma trovo  $\tilde{x}_h$

(cioè non è che lo risolvo ma gli faccio fare un po' di giri)



**EQ RESIDUO** • Calcolo  $r = b - A\tilde{x}_h$

**COARSE GRID** • "Proietta" su  $\Omega^{2h}$  COARSE GRID e risolvo  $A_{2h}\tilde{e}_{2h} = r_{2h}$

**CORRECTION** Poi "prolunga"  $\tilde{e}_{2h}$  a un  $\tilde{e}_h$  (per tornare alla griglia  $\Omega^h$ )

$$x_h = \tilde{x}_h + \tilde{e}_h$$

Dobbiamo trovare un modo per passare da  $\Omega^h$  a  $\Omega^{2h}$  e viceversa.

**(OPERATORE DI GRID TRANSFER)**

Chiamerò: **RESTRITTORE**  $r_h \rightarrow r_{2h}$

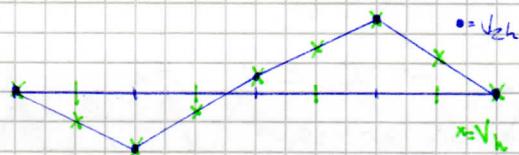
**PROLUNGATORE**  $\tilde{e}_{2h} \rightarrow \tilde{e}_h$

N.B.  $r_{2h}$  e  $\tilde{e}_{2h}$  non sono il residuo e l'errore "esatti" perché sono calcolati col restrittore ep...

Inoltre daremo considerazione cos'è  $A_{2h}$  (vedremo che come è intuitivo è la "restrizione" di  $A_h$  allo mesh più rado).

\* Vediamo prima come definiamo il PROLUNGATORE:  $I_{2h}^h: \Omega^{2h} \rightarrow \Omega^h$

$$I_{2h}^h v^{2h} = v^h \quad \text{t.c.} \quad \begin{cases} (v^h)_i = (v^{2h})_i & i = 0, \dots, N/2 \\ (v^h)_{2i+1} = \frac{1}{2}[(v^{2h})_i + (v^{2h})_{i+1}] & i = 0, \dots, N/2 - 1 \end{cases}$$



La matrice della trasformazione è:

$$C = \begin{bmatrix} \frac{1}{2} & & & & & & \\ & 1 & & & & & \\ & & \frac{1}{2} & & & & \\ & & & \frac{1}{2} & & & \\ & & & & 1 & & \\ & & & & & \frac{1}{2} & \\ & & & & & & \frac{1}{2} \end{bmatrix} \quad (7 \times 8 \text{ nodes del disegno})$$

\* **RESTRITTORE**:  $I_h^{2h}: \Omega^h \rightarrow \Omega^{2h}$  (limitazione)

$$I_h^{2h} v_h = v_{2h} \quad \text{t.c.} \quad (v^{2h})_i = (v^h)_{2i} \quad i = 0, \dots, N/2$$

Nel disegno sopra è come se parto dalla verde e poi tengo solo i punti blu. Quindi tengo inalterate le informazioni che mi servono.

$$\text{La matrice è } \begin{bmatrix} 0 & 1 & & & & & \\ & & 0 & 1 & & & \\ & & & & 0 & 1 & \\ & & & & & & 0 & 1 & \\ & & & & & & & & 0 & 1 & \\ & & & & & & & & & & 0 & 1 & \\ & & & & & & & & & & & & 0 & 1 \end{bmatrix} \quad (3 \times 7)$$

In realtà con questo restrittore la  $A_{2h}$  viene brutta quindi se ne usa un altro:

\* **RESTRITTORE**:  $I_h^{2h} v_h = v_{2h} \quad \text{t.c.} \quad (v^{2h})_i = \frac{1}{4}[(v^h)_{j-1} + 2(v^h)_j + (v^h)_{j+1}]$  **FULL-WEIGHTING**

$$\text{La matrice è } \begin{bmatrix} \frac{1}{4} & \frac{1}{2} & \frac{1}{4} & & & & \\ & \frac{1}{4} & \frac{1}{2} & \frac{1}{4} & & & \\ & & \frac{1}{4} & \frac{1}{2} & \frac{1}{4} & & \\ & & & \frac{1}{4} & \frac{1}{2} & \frac{1}{4} & \\ & & & & \frac{1}{4} & \frac{1}{2} & \frac{1}{4} \end{bmatrix} = \frac{1}{2} C^T$$

↳ cosa che ci piacerà molto

**SCHEMA di un PASSO di METODO:**

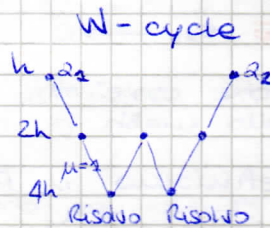
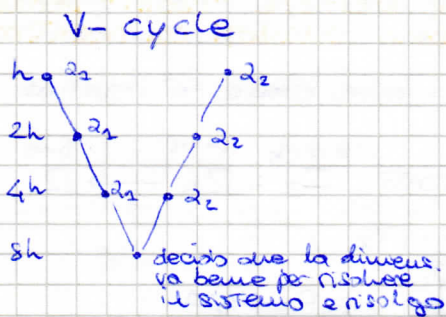
$$x_{\text{new}}^h = \text{Metodo}(x_{\text{old}}^h, b^h, A^h, a_1, a_2)$$



- > tipicamente  $2, 3$  non di più
- 1) Rilasso  $q_1$  volte  $A_h u_h = b_h$  in  $2^h$  con innesco arbitrario  $\tilde{x}_h^0$   
 $\leadsto \tilde{x}_h$  (PRESMOOTHING)
  - 2) Calcolo  $r_h = b_h - A_h \tilde{x}_h$
  - 3) Restrizione:  $\tilde{r}_{2h} = I_h^{2h} r_h$  proiezione del residuo (non è quello vero) in  $2^{2h}$
  - 4) "Risolvo"  $A_{2h} \tilde{e}_{2h} = \tilde{r}_{2h}$
  - 5) Prolungatore  $\tilde{e}_h = I_{2h}^h \tilde{e}_{2h}$
  - 6) Correggo approssimazione smoother:  $\hat{x}_h = \tilde{x}_h + \tilde{e}_h$
  - 7) Rilasso  $q_2$  volte  $A_h u_h = b_h$  (nel caso si fossero introdotte alte frequenze a causa degli errori nel non considerare il residuo esatto)  
 con innesco  $\hat{x}_h$  (POSTSMOOTHING)  
 $\leadsto x_h^{\text{new}}$

Questo metodo proposto così è detto TWO-GRID, se al punto (4) invece di risolvere perché la dim è troppo grande riapplica la ricorsione aggiungo una griglia...

Abbiamo più metodi per "approssimare":



19  
15/01/2013

## PROPRIETÀ DI CONVERGENZA

$\{Z_k\}_{k \geq 1}$  mesh t.c.  $Z_{k+1} \subset Z_k \rightarrow V_{k+1} \subset V_k \quad \forall k$  ( $Z_1$  mesh + coarse)  
 cioè la mesh successiva si ottiene per raffinamento.  
 Inoltre si suppone che la mesh non sia troppo brutta ("quasi-uniformi", "regolari")

## TEO: (di CONVERGENZA per il TWO GRID)

$\exists c > 0$  indipendente da  $k$  (indice di mesh) t.c.

$$\|\hat{e}_{k+1}\|_E \leq c m^{-1/2} \|e_0\|_E$$

dove abbiamo indicato con  $\|\cdot\|_E = (a(\cdot, \cdot))^{1/2}$  (norma energia) e

$$\hat{e}_{k+1} = \underbrace{\tilde{z}}_{\substack{\text{sol. esatta} \\ \text{FEM } A\tilde{z}=b}} - \underbrace{TG(\tilde{z}_0)}_{\substack{\text{innesco} \\ \text{\# passi di smoother}}}$$

Il teorema in pratica dice che il TG è una contrazione con fattore di contrazione indipendente da  $k$ .

## TEO: (di CONVERGENZA per V-CYCLE a $K$ livelli)

Sia  $m = \#$  passi di smoother, vale:

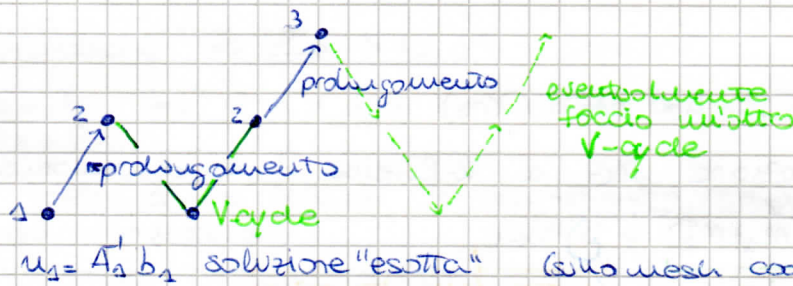
$$\|z - \underbrace{MG_m(K, z_0)}_{\substack{\text{soluzione esatta} \\ \text{FEM } A\tilde{z}=b}}\|_E \leq \underbrace{\frac{c^*}{c^* + m}}_{\substack{\text{livello di multigrid} \\ \text{approssimazione} \\ \text{data dal multigrid}}} \|z - z_0\|_E$$

Dunque  $\forall m$  MG è una contrazione e  $\chi$  è indipendente da  $K$  ( $= \#$  livelli di ricorsione)



# FULL MULTIGRID:

IDEA: Cercare l'innescio migliore  $\rightarrow$  così ottengo velocità convergenza maggiore



## ALGORITMO:

$k=1$   $\hat{u}_1 = A_1^{-1} b_1$  sol esatta (metodo diretto)

(più coarse)

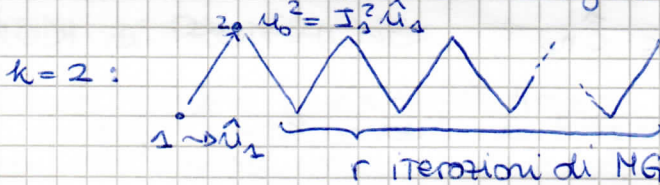
$k \geq 2$   $u_0^k = I_{k-1}^k \hat{u}_{k-1}$  innescio

$u_r^k = MG(k, u_{r-1}^k, b_k)$

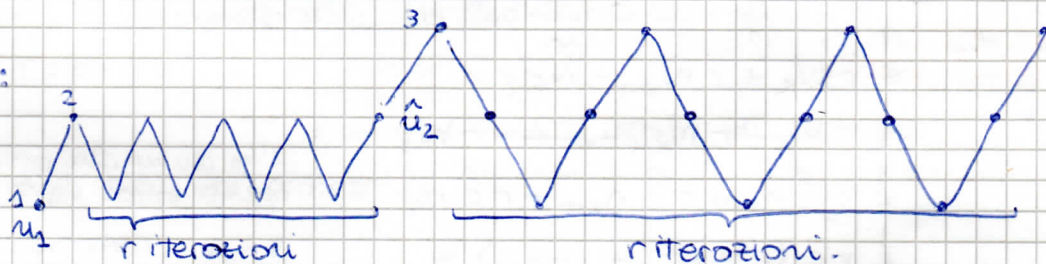
$1 \leq r \leq r$  [nel disegno ho preso  $r=1$ ]

$\hat{u}_k = u_r^k$

Facciamo di nuovo un disegno per copiare l'algoritmo



$k=3$ :



Ora vediamo un risultato che riesce a legare la soluzione approssimata alla soluzione nel continuo. (caso multibella, di solito legiamo l'approssimata alla soluzione esatta nel discreto)

$u$  = soluzione esatta del problema variazionale

$u_k \in V_k$  approssimazione a FE lineari di  $u$ .

Sapriamo (ved. Russo) che vale  $\|u - u_k\|_{H^s(\Omega)} \leq C h_k^{r-s} \|u\|_{H^r(\Omega)}$   $s=0,1$   $r=1,2$   $\rightarrow$  max diametro di  $T \in \mathcal{T}_k$

TEO: Se l'iterazione MG a  $k$  livelli è una contrazione con fattore indipendente da  $k$  allora  $\exists C > 0$  t.c.

$$\|u_k - \hat{u}_k\|_E \leq C h_k \|u\|_{H^2(\Omega)}$$

sol. app. FE  $\rightarrow$  appros. FULL MULTIGRID di  $u_k$

(è un risultato analogo a quello sopra  $\cap$ )

## COSTO COMPUTAZIONALE (LINEARE)

Per primo cosa vogliamo dare una stima asintotica di  $\dim V_k = n_k$

Per la formula di Eulero Abbiamo:  $V_k^I - e_k^I + t_k = 1$

vertici interni  $e_k^I$  edge int.  $t_k$  triangoli  $\mathcal{T}_k$

$$\Rightarrow n_k = V_k^I = 1 + e_k^I - t_k$$



Abbiamo che:  $\begin{cases} e_k^I = 2e_{k-1}^I + 3t_{k-1} \\ t_k = 4t_{k-1} \end{cases}$

da cui si ottiene:

$$n_k \sim \frac{t_1}{8} \cdot 4^k$$

**PROP:** le full multigrid ( $p=4$ ) ha costo  $O(n_k)$

dim: • Costo dell'iterazione MG a  $k$ -livello? Lo chiamiamo  $W_k$

$$W_k \leq \underbrace{C n_k (m_1 + m_2)}_{\substack{\text{costo prodotto} \\ \text{matrice vettore} \\ \text{MEMO: } A \text{ è sparsa}}} + \underbrace{p W_{k-1}}_{\substack{\text{costo delle iterazioni} \\ \rightarrow \# \text{ di diadrate: } p=1 \text{ Vcycle} \\ p=2 \text{ Wcycle}}}$$

$$\begin{aligned} &\leq C \tilde{m} n_k + p(C \tilde{m} n_{k-1} + p W_{k-2}) = \\ &\dots \leq C \tilde{m} n_k + p C \tilde{m} n_{k-1} + \dots + p^{k-1} C \tilde{m} n_1 = \\ &\stackrel{!}{=} \hat{C} \tilde{m} 4^k + p \hat{C} \tilde{m} 4^{k-1} + \dots + p^{k-1} \hat{C} \tilde{m} 4 = \quad \text{con } \hat{C} = \frac{C t_1}{8} \\ &\stackrel{!}{=} \hat{C} \tilde{m} 4^k \left( 1 + \frac{p}{4} + \dots + \frac{p^{k-1}}{4^{k-1}} \right) = \\ &\stackrel{se}{p/4 < 1} \leq \hat{C} \tilde{m} 4^k \frac{1}{1-p/4} \leq \hat{C} n_k \end{aligned}$$

•  $\hat{W}_k$  = costo full multigrid a livello  $k$

$$\begin{aligned} \hat{W}_k &= \hat{W}_{k-1} + \underbrace{r W_k}_{\substack{\text{\# iterazioni di multigrid}}} = \\ &\leq r W_k + r W_{k-1} + \hat{W}_{k-2} = \\ &\leq r W_k + r W_{k-1} + \dots + r W_1 = \\ &\leq r \hat{C} n_k + r \hat{C} n_{k-1} + \dots + r \hat{C} n_1 = \quad \text{l'ho chiamata così perché è quella sopra, poi le chiamo sempre C} \\ &\leq r C (4^k + 4^{k-1} + \dots + 4) = \\ &\leq r C 4^k \left( 1 + \frac{1}{4} + \dots + \frac{1}{4^{k-1}} \right) \leq r C 4^k \frac{1}{1-1/4} \leq C n_k \quad \blacksquare \end{aligned}$$

Ora torniamo indietro e cerchiamo di vedere nuovamente il metodo multigrid come metodo iterativo stazionario.

Esaminiamo il TWOGRID per semplicità. [Presmoother + coarse grid correction]

$$\begin{aligned} \tilde{x}^h &= G_h^{\text{pre}} x_{\text{old}}^h + c \\ \begin{cases} r^h = b^h - A x^h \\ \tilde{r}^{2h} = I_h^{2h} r^h \\ A_{2h} \tilde{e}^{2h} = \tilde{r}^{2h} \\ \tilde{e}^h = I_h^h \tilde{e}^{2h} \\ x_{\text{new}}^h = \tilde{x}^h - \tilde{e}^h \end{cases} \quad \text{CGC} \end{aligned}$$

$$x_{\text{new}}^h = \tilde{x}^h + I_h^h \tilde{e}^{2h} = \tilde{x}^h + I_h^h A_{2h}^{-1} \tilde{r}^{2h} = \dots = \overset{\text{identità}}{(I_h - I_h^h A_{2h}^{-1} I_h^{2h} A_h)} \tilde{x}^h + \dots$$

so uso post-smoother  
comporre  $G_h^{\text{post}}$   
qui davanti

$(I_h - I_h^h A_{2h}^{-1} I_h^{2h} A_h) G_h^{\text{pre}} x_{\text{old}}^h + \dots$   
C<sub>h</sub> coarse grid correction

⇒ METODO ITERATIVO  
STAZIONARIO



OSS: Dato che è un metodo iterativo stazionario la decrescita dell'errore può essere al più lineare (quindi non può comportarsi come il CG in presenza di cluster - od esempio -)

def:  $P$  **PROIETTORE (OBLIQUO)** se  $P^2 = P$

OSS: se  $P$  proiettore  $\Rightarrow \Sigma P \leq \{0, 1\}$  (spettro)

$$\text{dim: } Px = \lambda x \Rightarrow P^2 x = \lambda^2 x \Rightarrow \lambda^2 - \lambda = 0 \Rightarrow \lambda(1-\lambda) = 0$$

$$P^2 x = Px$$

Dunque  $\lambda(P) = 1$  e quindi  $P$  non è un metodo iterativo stazionario convergente.

**PROP:**  $C_H = I_n - I_n^H A_{2H}^{-1} I_n^{2H} A_H$  è un proiettore.

$$\text{dim: } C_H^2 = (I_n - I_n^H A_{2H}^{-1} I_n^{2H} A_H)(I_n - I_n^H A_{2H}^{-1} I_n^{2H} A_H) = I_n - 2 I_n^H A_{2H}^{-1} I_n^{2H} A_H + I_n^H A_{2H}^{-1} I_n^{2H} A_H I_n^H A_{2H}^{-1} I_n^{2H} A_H$$

$$= I_n - I_n^H A_{2H}^{-1} I_n^{2H} A_H = C_H$$

sono uguali

Si può però dimostrare che  $C_H$  è un proiettore ortogonale se si usa il prodotto scalare  $(\cdot, \cdot)_A$  (quello del CG) [Questo è un legame con CG ma non ho capito perché]

16/01/2013  
(20)

## METODI di PROIEZIONE

Vogliamo risolvere il sistema  $Ax = b$ . L'idea è quella di cercare una soluzione approssimante in un certo sottospazio  $K \subseteq \mathbb{R}^n$ , dim  $K = m$ . (ossia soddisfacente determinati vincoli).

Tipicamente si scelgono dei vincoli di ortogonalità:

$$r = b - Ax \perp m \text{ vettori e.c.}$$

CONDIZIONI di  
PETROV - GALERKIN

$$m \text{ vincoli in } L \text{ stsp dim } L = m$$

Ci sono due tipi di proiettori:

- ORTOGONALI  $L = K$
- OBLIQUI  $L \neq K$

Il problema è quindi:

$$\text{Trovare } \tilde{x} \in K \text{ sol. appross. t.c. } b - A\tilde{x} \perp L$$

che introducendo il vettore di innesco  $x_0$  possiamo riscrivere:

$$\text{Trovare } \tilde{x} \in x_0 + K \text{ t.c. } b - A\tilde{x} \perp L$$

→ Metodi di proiezione in stsp di Krylov.

def:  $K_m(M, v)$   $M \in \mathbb{C}^{n \times n}$   $v \in \mathbb{C}^n$ ,  $K_m(M, v) := \langle v, Mv, M^2 v, \dots, M^{m-1} v \rangle$   
 $= \{ u \in \mathbb{C}^n \text{ t.c. } u = \sum_{i=0}^{m-1} p_i(M) v \}$

def: **METODO di PROIEZIONE in STSP di KRYLOV** per risolvere  $Ax = b$ . Cerco  $x_k$  soluzione approssimante in  $x_0 + K_k$  con  $x_0$  vettore di innesco arbitrario, con

$$K_k = K_k(A, r_0) = \langle r_0, Ar_0, \dots, A^{k-1} r_0 \rangle$$

e inoltre si impone  $r_k = b - Ax_k \perp L_k \text{ stsp t.c. dim } L_k = k$ .

Abbiamo ancora diverse possibilità di scelta perché dobbiamo fissare  $L$ .

$$\begin{aligned} &\nearrow L_k = K_k(A, r_0) \text{ o } L_k = AK_k(A, r_0) \text{ (GMRES)} \\ &\nwarrow L_k = K_k(A^T, r_0) \end{aligned}$$



1° PASSO: Costruire una base ortonormale  $K_k(A, q_0)$ .

Se abbiamo  $v_1 \rightarrow v_k$  vettori linearmente indipendenti possiamo ortonormalizzarli col metodo di GRAM-SCHMIDT (classico).

Metodo classico **GRAM-SCHMIDT**:

$v_1 \rightarrow v_n$  e.i.

$$\tilde{u}_k = v_k - \sum_{i=1}^{k-1} \langle v_k, u_i \rangle u_i$$

$$u_k = \tilde{u}_k / \|\tilde{u}_k\|_2$$

Metodo **GM modificato** (per l'aritmetico floating point)

$$u_1 = v_1 / \|v_1\|_2$$

per  $k=2 \rightarrow n$

$$\tilde{u}_k = v_k$$

per  $i=1 \rightarrow k-1$

$$\tilde{u}_k = \tilde{u}_k - \langle \tilde{u}_k, u_i \rangle u_i$$

$$u_k = \tilde{u}_k / \|\tilde{u}_k\|_2$$

procedo a tope rispetto al vettore  
già modificato.

(È meglio in aritmetica floating point  
perché è come se avessi  $\tilde{u}$  più pulito)

Metodo di **ARNOLDI** = Metodo GM modificato per costruire base ortonormale per  $K_k(A, q_0)$

INPUT:

?  $A, q_0$  ?

OUTPUT:  $q_1 \rightarrow q_k$  vettori ortonormali e base per  $K_k(A, q_0)$

$q_1$  etc.  $\|q_1\|_2 = 1$  (abbiamo già lo normalizzato)

per  $j=1 \rightarrow k$

$$\tilde{q}_{j+1} = A q_j$$

per  $i=1 \rightarrow j$  ← è giusto perché è già avanti di un passo

$$h_{ij} = \langle \tilde{q}_{j+1}, q_i \rangle$$

$$\tilde{q}_{j+1} = \tilde{q}_{j+1} - h_{ij} q_i$$

$$h_{j+1,j} = \|\tilde{q}_{j+1}\|_2$$

$$q_{j+1} = \tilde{q}_{j+1} / h_{j+1,j} \quad (\bullet)$$

**OSS:** Non è ovvio che quella che abbiamo trovato sia una base per  $K_k(A, q_0)$   
(è ovvio che sono ortonormali ma non che sia base)

**PROP:**  $\{q_1 \rightarrow q_k\}$  dell'algoritmo di ARNOLDI formano una base ortonormale per  $K_k(A, q_0)$

dim: Sono ortonormali per costruzione. Vogliamo verificare che

$$q_j = P_{j-1}(A) q_0$$

basta che scelga:

Per induzione:  $j=1$  vero  $q_1 = p_0(A) q_0$

$$p_0(t) \equiv 1$$

hp di induzione  $i \leq j$   $q_i = p_{i-1}(A) q_0$

Verifichiamo per  $i=j+1$ . Per def dell'algoritmo abbiamo:

$$h_{j+1,j} q_{j+1} = \tilde{q}_{j+1} = A q_j - \sum_{i=1}^j h_{ij} q_i =$$



$$= A p_{j-1}(A) q_1 - \sum_{i=1}^j h_{ij} p_{i-1}(A) q_1 = \tilde{p}_j(A) q_1$$

PROP: L'algoritmo di Arnoldi si può riscrivere come:

$$A Q_k = Q_k H_k + h_{k+1,k} \overset{\text{vettore}}{q_{k+1}} e_k^T = Q_{k+1} H_{k+1,k}$$

HOUSEHOLDER-ARNOLDI

$$\begin{matrix} \boxed{A}^n & \boxed{Q_k}^n & = & \boxed{Q_k}^n & \boxed{H_k}^k & + & \boxed{\begin{matrix} \circ & \vdots \\ \vdots & \end{matrix}} & = & \boxed{Q_{k+1}}^{k+1} & \boxed{H_{k+1,k}}^{k+1} \end{matrix}$$

ovv.  $Q_k = [q_1 | \dots | q_k]$   $H_k = \begin{bmatrix} h_{11} & & \\ & \ddots & \\ 0 & & h_{kk} \end{bmatrix}$  Hessenberg superiore

dim:  $j=1$  per  $i=1$   $\tilde{q}_2 = A q_1 - h_{11} q_1$

$$\rightarrow A q_1 = h_{11} q_1 + \tilde{q}_2 = h_{11} q_1 + h_{21} q_2$$

$j=2$  per  $i=1$   $\tilde{q}_3 = A q_2 - h_{12} q_1$

$i=2$   $\tilde{q}_3 = \tilde{q}_3 - h_{22} q_2 = A q_2 - h_{12} q_1 - h_{22} q_2$  (vedi l'algoritmo per verificare le uguaglianze che sono giuste, conviene ->)

$$\rightarrow A q_2 = h_{12} q_1 + h_{22} q_2 + \tilde{q}_3 = h_{12} q_1 + h_{22} q_2 + h_{32} q_3$$

Quindi in generale ricaviamo:

$$A q_k = h_{1k} q_1 + h_{2k} q_2 + \dots + h_{kk} q_k + h_{k+1,k} q_{k+1}$$

GMRES (matrici non Hermitiane)

E' un metodo di proiezione su:

$$K_k(A, q_1) = \langle q_1, A q_1, \dots, A^{k-1} q_1 \rangle$$

$$q_1 = \frac{r_0}{\|r_0\|_2}$$

$$L_k = A K_k(A, q_1)$$

Cerca  $x_k = x_0 + K_k(A, q_1)$

ne ho una base ortonormale  $q_1, \dots, q_k$   
 $\rightarrow$  NB e' un modo per scrivere una comb. lineare dei vettori  $q_i$  (con  $\gamma_i$  come coefficienti)

$$= x_0 + Q_k y$$

il mio obiettivo e' minimizzare la norma 2 del residuo:  $\|r_k\|_2$

$$r_k = b - A x_k = b - A(x_0 + Q_k y) = r_0 - A Q_k y$$

Poiché  $\beta = \|r_0\|_2$   $e_1 = [1, 0, \dots, 0]^T \in \mathbb{C}^{k+1}$  per costruzione di Arnoldi abbiamo:

$$Q_{k+1} e_1 = q_1 := \frac{r_0}{\beta}$$

E quindi abbiamo  $r_0 = \beta Q_{k+1} e_1$  (M)

Pb: Trovare  $\min_{y \in \mathbb{C}^k} \|r_k\|_2$ . Abbiamo

$$\begin{aligned} \min_{y \in \mathbb{C}^k} \|r_k\|_2 &= \min_{y \in \mathbb{C}^k} \|r_0 - A Q_k y\|_2 \\ &= \min_{y \in \mathbb{C}^k} \|\beta Q_{k+1} e_1 - Q_{k+1} H_{k+1,k} y\|_2 \\ &\stackrel{+PROP}{=} \min_{y \in \mathbb{C}^k} \|Q_{k+1} (\beta e_1 - H_{k+1,k} y)\|_2 \\ &= \min_{y \in \mathbb{C}^k} \|\beta e_1 - H_{k+1,k} y\|_2 \quad (B) \end{aligned}$$

$$MEMO: Q_{k+1}^H Q_{k+1} = I$$

La matrice ha rango massimo per costruzione  
 $\Rightarrow$  Pb lineare di minimi quadrati ha un'unica soluzione.



**TEO:** Dato il problema ai minimi quadrati lineari:

(MEMO) "Determinare  $X = \{x \in \mathbb{C}^m \mid \|Ax - b\|_2 = \min_{y \in \mathbb{C}^n} \|Ay - b\|_2 \quad A \in \mathbb{C}^{m \times n}\}$ "

Vali:

- $x \in X \Leftrightarrow A^H A x = A^H b$
- $X$  è insieme non vuoto, chiuso e connesso
- $X = \{x^*\}$  sse  $\text{rk}(A)$  massimo
- $\exists x^* \in X$  t.c.  $\|x^*\|_2 = \min_{x \in X} \|x\|_2$

Quindi in particolare nel nostro caso dato che il rango è max abbiamo un'unica soluzione (che quindi è anche il minimo!).

Sia  $H_{k+1,k} = F^H R$  fattorizzazione QR di una matrice rettangolare (F sarebbe la Q ma ho già troppe Q in giro)

$$\begin{array}{c} k+1 \\ \boxed{H_k} \\ \begin{array}{cc} 0 & h_{kk} \\ \hline & \end{array} \\ k \end{array} = \begin{array}{c} k+1 \\ \boxed{F^H} \\ k \end{array} \begin{array}{c} k \\ \boxed{R} \\ \begin{array}{cc} & \\ 0 & \\ \hline & D \end{array} \\ k \end{array} \left. \begin{array}{l} \\ \end{array} \right\} \text{tridiag sup.}$$

Dunque:

$$\begin{aligned} (*) &= \min_{y \in \mathbb{C}^k} \|\beta e_1 - F^H R y\|_2 = \\ &= \min_{y \in \mathbb{C}^k} \|F^H (\beta F e_1 - R y)\|_2 = \text{altro siccome l'ultima riga è nulla considero solo le prime } k \text{ componenti, e minimizzo su quelle, che è uguale,} \\ &= \min_{y \in \mathbb{C}^k} \|\beta F e_1 - R y\|_2 = \min_{y \in \mathbb{C}^k} \|(\beta F e_1)_{k+1} - R_{kk} y_{k+1}\| \end{aligned}$$

Ora abbiamo ottenuto un sistema quadrato con  $R_{kk}$  non singolare!  
Non dobbiamo più risolvere un problema di minimizzazione!

$\Rightarrow y \in \mathbb{C}^k$  cercato è la soluzione di:

$$R_{kk} y = \beta (F e_1)_{k+1}$$

e vale inoltre  $\|r_k\|_2 = \|b - Ax_k\|_2 = \|\beta F e_1 - R y\|_2 = \beta (F e_1)_{k+1}$   
mi rimane solo l'ultima componente perché sulle prime  $k =$

**OSS:** Questo metodo è abbastanza costoso! Ma d'altronde i sistemi non hermitiani sono cose brutte (non sono animali domestici :-))  
Costano soprattutto in termini di memoria.

Ci sono altre versioni del GMRES oltre a quella "classica":

- con RESTART
  - con TRONCAMENTO
- } derivano dal problema del costo di memoria

Idea restart: Faccio m iterazioni  $q_1, \dots, q_m$  mi fermo e prendo l' $x_m$  trovato come nuovo innescio!

Idea truncamento: Non faccio l'ortogonalizzazione completa, cioè non torno indietro fino all'indice 1. Faccio una "quasi-ortogonalizzazione".

**ALGORITMO:** (per GMRES)

[Saad p 158]

1. Calcolo  $r_0 = b - Ax_0$ ,  $\beta = \|r_0\|_2$ ,  $v_1 = r_0/\beta$   
(alloca la memoria)
2. Dimensionare  $H_m = [h_{ij}]_{i=1 \dots m+1; j=1 \dots m} = 0$
3. per  $j = 1 \dots m$   
Calcolo  $w_j = A v_j$

N.B. l'ha scritto per forza notore che usiamo tanta memoria

21  
18/01/2013



$$\left[ \begin{array}{l} \text{per } i=1 \rightarrow j \\ h_{ij} = (w_j, v_i) = w_j^T v_i \\ w_j = w_j - h_{ij} v_i \\ h_{j+1,j} = \|w_j\|_2 \\ v_{j+1} = w_j / h_{j+1,j} \end{array} \right.$$

4. Calcolare  $y_m$  minimizzatore di  $\| \beta e_1 - \tilde{H}_m y \|_2$  [QR e risolve " $Ry = b$ "]  
e  $x_m = x_0 + V_m y_m$  ove  $V_m = [v_1 | v_2 | \dots | v_m]$ .

### ANALISI PROPRIETÀ di CONVERGENZA (GMRES)

Al passo  $k$  ho:  $x_k = x_0 + z_k$  con  $z_k \in K_k(A, r_0) = \langle r_0, A r_0, \dots, A^{k-1} r_0 \rangle$

$$r_k = b - A x_k = b - A x_0 - A z_k = r_0 - A z_k \in \langle r_0, A r_0, \dots, A^{k-1} r_0 \rangle$$

Ossia:  $r_k = p_k(A) r_0$  con  $p_k$  polinomio di grado  $k$  t.c.  $p_k(0) = 1$   
$$= r_0 + \sum_{i=1}^k \gamma_i A^i r_0$$

Dunque possiamo scrivere la norma dell'errore (lo minimizza perché  $r_0$ )

$$\|r_k\|_2 = \min_{\substack{p_k \in \mathcal{P}_k \\ p_k(0)=1}} \|p_k(A) r_0\|_2$$

Ora vogliamo trovare una maggioranza per il residuo.

**PROP:** Se  $A = V \Lambda V^{-1}$  con  $V$  matrice con colonne gli autovettori (normalizzati in modo opportuno)

$$\Rightarrow \frac{\|r_k\|_2}{\|r_0\|_2} \leq \underbrace{K_2(V)}_{\text{NEW ENTRY!}} \cdot \min_{\substack{p_k \in \mathcal{P}_k \\ p_k(0)=1}} \left\{ \max_{i=1 \rightarrow n} |p_k(\lambda_i(A))| \right\}$$

**OSS:** Se  $A$  è normale  $\Rightarrow V$  unitaria e dunque  $\|V\|_2 = \|V^{-1}\|_2 = 1$   
 $\Rightarrow K_2(V) = 1$  e il termine "malo" sparisce.

Se  $A$  non è normale in generale  $K_2(V) \neq 1$  e ce lo dobbiamo tenere!

$$\text{d'ui: } \|r_k\|_2 = \min_{\substack{p_k \in \mathcal{P}_k \\ p_k(0)=1}} \|p_k(A) r_0\|_2 = \min_{\substack{p_k \in \mathcal{P}_k \\ p_k(0)=1}} \|p_k(V \Lambda V^{-1}) r_0\|_2 =$$

$$\text{dato che } A^s = V \Lambda^s V^{-1} \leftarrow \begin{array}{l} p_k \in \mathcal{P}_k \\ p_k(0)=1 \end{array}$$

$$= \min_{\substack{p_k \in \mathcal{P}_k \\ p_k(0)=1}} \|V p_k(\Lambda) V^{-1} r_0\|_2 \leq \min_{\substack{p_k \in \mathcal{P}_k \\ p_k(0)=1}} \|V\|_2 \|p_k(\Lambda)\|_2 \|V^{-1}\|_2 \|r_0\|_2 =$$

$$\stackrel{\text{def}}{=} \min_{\substack{p_k \in \mathcal{P}_k \\ p_k(0)=1}} \|p_k(\Lambda)\|_2 \|r_0\|_2 \cdot K_2(V)$$

$$\stackrel{\text{d'ui}}{=} \frac{1}{K_2(V)} \min_{\substack{p_k \in \mathcal{P}_k \\ p_k(0)=1}} \left\{ \max_{i=1 \rightarrow n} |p_k(\lambda_i(A))| \right\} \|r_0\|_2$$

**OSS<sub>1</sub>:** Non è noto se questa è una stima stretta! (diversamente che per il GC), questo perché non è detto che il polinomio che minimizza con dentro  $r_0$  sia lo stesso che minimizza senza  $r_0$ .

Se  $A$  è normale però ho  $K_2(V) = 1$  e dunque la limitazione è stretta.

**OSS<sub>2</sub>:** Anche questo metodo è sensibile alla presenza di cluster. Se i cluster sono lontani dall'origine si comporta bene, se sono intorno all'origine no.

**OSS<sub>3</sub>:** Se  $A$  non è normale, ma  $V$  è "abbastanza" ben condizionata allora la PROP ci dà una stima non stretta ma ragionevole. Inoltre sono sostanzialmente solo gli autovettori a determinare le proprietà di convergenza (il ruolo degli autovettori è marginale).



Un altro strumento che viene usato per studiare (le proprietà) la velocità di convergenza è il seguente:

$$\text{Si consideri } F(A) = \{y^H A y \mid y \in \mathbb{C}^n, y^H y = 1\} \quad (\text{FIELD OF VALUES}), A \in \mathbb{C}^{n \times n} \\ = \left\{ \frac{y^H A y}{y^H y} \mid y \in \mathbb{C}^n, y \neq 0 \right\}$$

Si nota che questo caso "assomiglia" ai quozienti di Rayley. Quindi l'idea è quella di partire da quello che sappiamo nel caso simmetrico e cercare di estenderlo al caso non simmetrico.

**PROBLEMI DI CONVEZIONE - DIFFUSIONE - FEM**  
(streamline diffusion, SUPG = streamline upwind Petrov-Galerkin)

$$\begin{cases} -\operatorname{div}(\varepsilon \nabla u) + \beta \nabla u = f & \text{su } \Omega \\ u = 0 & \text{su } \partial\Omega \end{cases}$$

FEM:  $a(u, v) = F(v) \quad \forall v \in V$

$$a(u, v) = \int_{\Omega} \varepsilon \nabla u \cdot \nabla v + \int_{\Omega} \beta \cdot \nabla u \, v$$

$$F(v) = \int_{\Omega} f v$$

def:  $Pe$  **NUMERO DI PECLLET** che misura rapporto tra diffusione e convezione.

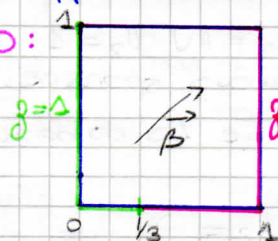
Si può mostrare che se  $Pe > 1$  localmente, FEM non funziona. Petro meglio.

Possiamo definire anche un numero di Peclet locale.

$$Pe_k = \frac{\|\vec{\beta}_k\| h_k}{2\varepsilon_k} \quad h_k \text{ diam } K \text{ con } K \in \mathcal{T} \text{ triangolo.}$$

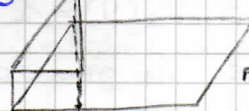
Se  $Pe_k > 1$  FEM non funziona, ma proprio non funziona il metodo dell'approssimazione con gli elementi finiti!

**ESEMPIO:**



$$\vec{\beta} = \begin{bmatrix} 1 \\ 3 \end{bmatrix}$$

( $g$  = dato di bordo)



$\approx$  c'è un salto nelle condizioni di bordo.

Noi ora fisseremo  $h_k$  e  $\vec{\beta}$  e vedremo cosa succede facendo tendere  $\varepsilon \rightarrow 0$ . (negli esempi, credo)

Vediamo uno schema del metodo SUPG. Consideriamo il problema (già proiettato sullo spazio più piccolo - triangoli)

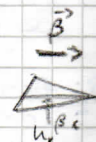
$$\int_{\Omega} \varepsilon \nabla u_L \cdot \nabla v_L + \int_{\Omega} (\vec{\beta} \nabla u_L) \cdot v_L + \sum_{K \in \mathcal{T}} \tau_K \int_K (\vec{\beta}_K \nabla u_L) (\vec{\beta}_K \nabla v_L) =$$

$$= \int_{\Omega} f v_L + \sum_{K \in \mathcal{T}} \tau_K \int_K f_K (\vec{\beta}_K \nabla v_L) \quad \rightarrow \text{ho correzione del termine noto e anche il termine ass. coef. della matrice}$$

ora abbiamo posto:

$$\tau_K = \frac{h_K^{\beta_K}}{2\|\vec{\beta}_K\|}$$

con  $h_K^{\beta_K}$  è il "diametro nella direzione di  $\vec{\beta}$ "



se  $Pe_k > 1$  e altrimenti  $\tau_k = 0$  (se  $Pe_k < 1$ ).

Vediamo come possiamo approssimare questi vari integrali che abbiamo introdotto.



$$\sum_k \int_K \vec{\beta} \cdot \nabla \varphi_k \nabla \varphi_j \stackrel{\text{formula del baricentro}}{\cong} \sum_k \frac{\text{dot}(\vec{\beta}_k(b_k), e_j^\perp) \cdot \text{dot}(\vec{\beta}_k(b_k), e_i^\perp)}{4|\Gamma|}$$

$$\sum_k \int_K f_k (\vec{\beta}_k \cdot \nabla \varphi_i) \cong \sum_k \frac{f_k(b_k)}{2} \text{dot}(\vec{\beta}_k(b_k), e_i^\perp)$$

**ESEMPIO:**  $g=0$  (dirichlet) trovare un polinomio che è un coseno  
 (da provare con matlab)  $f=0$   
 $\vec{\beta}$  due ruota intorno al punto  $(1,0) \Rightarrow \vec{\beta} = \begin{bmatrix} y-1 \\ -(x-0) \end{bmatrix}$

